

Conditional co-occurrence probability acts like frequency in predicting fixation durations

James K. Y. Ong

Upper Austria University of Applied Sciences

Reinhold Kliegl

University of Potsdam

The predictability of an upcoming word has been found to be a useful predictor in eye movement research, but is expensive to collect and subjective in nature. It would be desirable to have other predictors that are easier to collect and objective in nature if these predictors were capable of capturing the information stored in predictability. This paper contributes to this discussion by testing a possible predictor: conditional co-occurrence probability. This measure is a simple statistical representation of the relatedness of the current word to its context, based only on word co-occurrence patterns in data taken from the Internet. In the regression analyses, conditional co-occurrence probability acts like lexical frequency in predicting fixation durations, and its addition does not greatly improve the model fits. We conclude that readers do not seem to use the information contained within conditional co-occurrence probability during reading for meaning, and that similar simple measures of semantic relatedness are unlikely to be able to replace predictability as a predictor for fixation durations. **Keywords:** Co-occurrence probability, Cloze predictability, frequency, eye movement, fixation duration

Keywords: Co-occurrence probability, Cloze predictability, frequency, eye movement, fixation duration

Introduction

The subjective predictability of a word given a previous partial sentence context (specifically, *Cloze predictability*, see Taylor, 1953) has proven to be a useful predictor of the fixation duration on that word during reading (Kliegl, Nuthmann, & Engbert, 2006; Rayner, Ashby, Pollatsek, & Reichle, 2004). Unfortunately, predictability estimates are expensive to collect. A number of alternative measures have been used to represent aspects of word predictability, including transitional probability (McDonald & Shillcock, 2003; Frisson, Rayner, & Pickering, 2005), surprisal (Boston, Hale, Kliegl, Patil, & Vasishth, 2008), orthographic familiarity or regularity (White, 2008; White & Liversedge, 2006), semantic constraint (Pynte, New, & Kennedy, in press), and global sentence properties (Pynte & Kennedy, 2006); however, none of these measures have been able to displace predictability as a predictor of fixation durations. In this paper, we test another simple statistical measure derived from the Internet, the *conditional co-occurrence probability* (CCP), which quantifies the chance that a word occurs given its preceding context.

Motivation for choosing CCP

Linguistic measures derived from the Internet are interesting for researchers for two main reasons: they are easy to collect, and they contain nontrivial information about statistical regularities.

Turney (2001) compared the performance of two methods to correctly recognise synonyms, one based on a pointwise mutual information measure derived from the Internet, and the other derived from Latent Semantic Analysis (Landauer & Dumais, 1997). He found that the Internet-based method does at least as well as the other method and concludes that the amount of Internet training data compensates for the simplicity of the mutual information measure.

In addition, Keller and Lapata (2003) showed that frequency counts of word bigrams (pairs of neighbouring words) generated from the Internet can better mirror human plausibility judgements than the equivalent counts derived from clean corpora. The reduction in data sparsity due to the far larger source data seems to outweigh the inherent noisiness of web data. Their results also show that CCP performs better in a pseudo-disambiguation task than joint probability.

Based on these results, CCP presents itself as a candidate to capture part of the semantic relationship between a word and its context, and thus also as a possible partial replacement for Cloze predictability in predicting fixation durations during reading. In addition, eye movements during reading also point to the potential relevance of non-local semantic measures like CCP for predictability modelling. The accuracy of long regressive saccades (Kennedy & Murray, 1987) is evidence for the retention of the meaning of previously read words, while the very presence of long regressive saccades occur implies that previous words continue to be relevant long after they have been initially fixated.

In the rest of the paper, we will take a closer look at CCP and discuss its relevance to reading eye movements, as deduced from a repeated measures regression analysis.

Method

Definition of CCP

The probability of word w_j occurring, given that word w_i occurs in the context, $p(w_j|w_i)$, which we will call the conditional co-occurrence probability (CCP), can be calculated as follows:

$$p(w_j|w_i) = \frac{p(w_i \cap w_j)}{p(w_i)} = \frac{f(w_i \cap w_j)}{f(w_i)}.$$

This means that one takes a collection of “documents”, counts the number of documents that contain both the target word (w_j) and the context word (w_i), and divides this total by the number of documents containing the context word.

In our case, the frequency counts for words and co-occurring pairs of words were collected from the Google, Yahoo! and MSN search engines via their Application Programming Interfaces in August 2006. All word unigram frequencies were above zero, and zero word co-occurrence frequencies were increased to one to allow for log transformation of the resulting CCP.

The measure of relatedness between a word and its previous context C was taken to be the maximum of the CCPs of the target word paired with all other words in the context:

$$p(w_j|C) = \max_{w_i \in C} p(w_j|w_i).$$

Later, in the regression analyses, we will denote $\log p(w_j|C)$ as q_j .

Comparison of conditional co-occurrence with other lexical features

Figure 1 shows the relationship between CCP and Cloze predictability estimates for the German *Potsdam Sentence Corpus* (described in Kliegl, Grabner, Rolfs, & Engbert, 2004). Conditional co-occurrence probability and predictability are not unrelated, which is reflected in a correlation coefficient of about 0.5 (regardless of the choice of search engine) for the words in the Potsdam Sentence Corpus. However, it is clear that conditional co-occurrence only captures a small part of the information contained within predictability. Very low conditional co-occurrence probabilities tend to imply low predictability; this means that high CCP is a necessary but not sufficient condition for high predictability.

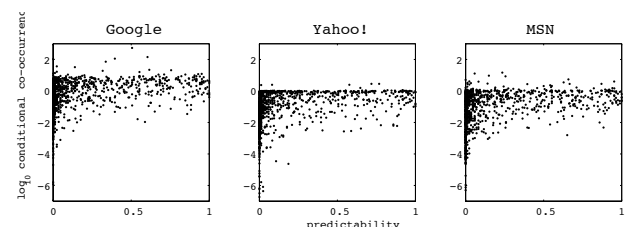


Figure 1. Conditional co-occurrence probability plotted against predictability for words in the Potsdam Sentence Corpus. Words that only rarely co-occur with their context tend to be unpredictable. Note that some of the results show conditional co-occurrence probabilities greater than one (or equivalently, the logarithm greater than zero); this anomaly occurs because we directly use the frequency estimates produced by the search engines, which exhibit both noise and systematic errors.

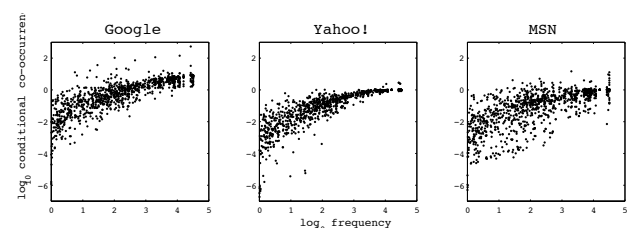


Figure 2. Conditional co-occurrence probability plotted against frequency for words in the Potsdam Sentence Corpus. There is a marked correlation between the two parameters, especially in the results produced by the Google and Yahoo! search engines. Note that some of the results show conditional co-occurrence probabilities greater than one (or equivalently, the logarithm greater than zero); this anomaly occurs because we directly use the frequency estimates produced by the search engines, which exhibit both noise and systematic errors.

Figure 2 shows the relationship between CCP and frequency. The relationship between CCP and frequency is far closer, as shown by a high correlation coefficient (0.8 for Google and Yahoo!, and 0.7 for MSN) for the words in the Potsdam Sentence Corpus.

Regression analysis

To deduce the effects of lexical features of words on fixation durations, Kliegl et al. (2006) previously conducted a comprehensive repeated measures regression analysis. Here, we build on this analysis by adding the CCPs of the current and neighbouring words to the list of predictors. Two aspects of these new predictors are of interest: whether they account for variance in the data previously accounted for by other predictors, and whether they can account for extra unique variance. We analyse the utility of these predictors in explaining single fixation durations (fixation durations on words fixated only once in first-pass reading), and gaze durations (the total time that a word is fixated when it is first encountered in first-pass reading).

In the rest of the paper, we will refer to the currently fixated word as word n , the neighbour to the left as word $n - 1$, and the neighbour to the right as word $n + 1$. Effects originating from word n are called *immediacy* effects; those from word $n - 1$, *lag* effects; and those from word $n + 1$, *successor* effects.

The eye movement data used in this analysis are the same as those used by Kliegl (2007)¹, and our reference regression models, which we will call the *baseline models*, are updated versions of his final linear mixed-effects models (one for single fixation durations and one for gaze durations). Thus, we adopt the same encoding of the model parameters:

- *Fixation duration* is transformed into log fixation duration;
- *Saccade amplitude* is measured in letters;
- *Fixation position* in letters is divided by word length to give a relative fixation position;
- *Length* is transformed into reciprocal word length;
- *Frequency* is transformed into log frequency, and then nested within lexical status to produce two variables, which we will call content word frequency and function word frequency;
- *Predictability* is transformed into logit predictability;
- *Lexical status* is a binary variable that is zero for content words and one for function words;
- *Skipping status*, which is only defined for words $n - 1$ and $n + 1$, is a binary variable that is zero for fixated words and one for skipped words.

All non-binary variables were centred about zero. The regression analyses were performed with the *lmer* program (Bates, 2007) in the R environment (R Development Core Team, 2008). In the results section, we represent these encodings with the following variables: l_k (length), g_k (content word frequency), h_k (function word frequency), p_k (predictability), x_k (lexical status) and s_k (skipping status), where the subscript k specifies which word is meant.

To the baseline models, we add three new predictors: q_n , q_{n-1} , and q_{n+1} , where q_j is as defined before. Since q_{n-1} is technically not defined when $n = 1$, we remove fixations on the first word of the sentence from our analysis. After adding the new predictors, significant two-way and three-way interactions are also added to the model, following the procedure used by Kliegl (2007) to generate the baseline model. Pairwise interactions between predictors on word $n - 1$ and word $n + 1$ were not considered in this analysis. Significant interactions are added in the following order:

1. Pairwise interactions with word length;
2. Pairwise interactions with word frequency, nested within lexical status;
3. Pairwise interactions with predictability;
4. Interactions with the lexical status of words n , $n - 1$, and $n + 1$, and the skipping status of words $n - 1$ and $n + 1$.

After the significant interactions are added, simple *random effects* are included; these random effects allow each subject to have, for example, a subject-specific word length effect instead of just the average word length effect of all subjects. No predictors or interactions present in the baseline model are removed from the model, even when they fail to remain significant predictors after addition of the new terms. We will refer to these final models as the *expanded models*.

Results

In this section, we give a brief summary of the main effects found in the baseline models in first-pass reading, and then describe the effects on the models of adding the CCP predictors.

¹ Specifically, the data used are the right eye fixations. It is possible that the left eye may be fixating different words and thus contributing to lag or successor effects, but Kliegl et al. (2006) reported that their regression models remained mostly unchanged after a restriction to binocularly registered fixations, with the only differences being in the fixated within-word letter positions.

Single fixation durations

In this baseline model, many of the lexical features of words n , $n - 1$, and $n + 1$ have a significant effect on the fixation duration on word n . The signs of the main effects of the lexical parameters (linear terms only) are summarised in Table 1. Because we have restricted this particular analysis to single fixation cases, a number of the immediacy effects, like length or lexical status, have no significant effect on the fixation duration; the effects of these parameters do appear when multiple fixation cases are investigated (for comparison, see the analysis of gaze durations below).

Table 1
Main Effects of Lexical Parameters on Single Fixation Durations

Parameter	Word $n-1$	Word n	Word $n+1$
l	↗	.	.
g	↘	.	.
h	↘	↘	.
p	↘	↘	.
x	↘	.	↘
<i>new: q</i>	↘	.	.

Note. Only the effects of the linear terms are shown here. If increasing the parameter increases fixation duration, then a '↗' is displayed, or if the opposite is true, then a '↘' is displayed. If there is no statistically significant effect (at the 99.7% confidence level), then a '.' is shown. For example, when l_{n-1} is smaller (meaning that the previous word is longer, since l represents reciprocal length), the fixation duration on the current word is shorter.

We add the three predictors q_n , q_{n-1} , and q_{n+1} to the baseline model, and then successively add interactions as detailed in the results section. Addition of the random effects representing intersubject variability to both the baseline and the expanded models does not change the pattern of differences between the models significantly. For us then, the random effects only play a role in the amount of total variance explained by the model, which becomes relevant when we report log likelihood values later.

Only one of the three additional predictors, q_{n-1} , ends up having a significant main effect on fixation duration, and this turns out to be negative. This means that if word $n - 1$ co-occurs frequently with one of its preceding words, the fixation duration on word n is likely to be shorter. The inclusion of the new predictors causes corre-

sponding compensatory changes (in the expected directions) in the main effects of the frequency predictors, especially g_{n-1} and h_{n-1} ; these changes occur because of the high covariances between the CCP and frequency predictors (see Table 2). However, there is no significant change in the main effects of the predictability predictors.

Table 2
Covariance Matrix for the Single Fixation Durations Model

	g_{n-1}	g_{n+1}	p_n	p_{n-1}	p_{n+1}	x_n	x_{n-1}	x_{n+1}	q_n	q_{n-1}	q_{n+1}
g_n	↘	.	.
g_{n-1}		↘	.
g_{n+1}			↗	.	.	↘
p_n			
p_{n-1}				
p_{n+1}					
x_n						
x_{n-1}							
x_{n+1}									.	.	↘
q_n										.	.
q_{n-1}											.

Note. The predictors shown are content word frequency (g), predictability (p), lexical status (x) and CCP (q), for words n , $n - 1$, and $n + 1$. A positive covariance more positive than 0.25 is represented by a '↗', and a negative covariance more negative than -0.25 is represented by a '↘'. Double arrows signify covariances with magnitudes larger than 0.50. Small covariances (magnitude less than 0.25) are represented by a '.'.

Most of the interactions from the baseline model remain unchanged after the addition of interactions with the new predictors. Notable exceptions are $x_n : x_{n-1}$, $x_n : g_{n-1}$, and $x_n : h_{n-1}$, which change in the expected directions to compensate for the addition of the new interactions; one example is the effect of the interaction $x_n : g_{n-1}$, which is highly nonindependent from the interaction $x_n : q_{n-1}$. Interactions in the baseline model containing predictability do not change significantly.

The inclusion of the new predictors causes a small improvement in the fit of the model; the increase in log likelihood is shown in Table 3. However, this improvement is counterbalanced by the increase in complexity of the model, which is reflected in the resulting *increase* in the Bayesian Information Criterion. The removal of non-significant predictors left over from the baseline model helps to improve the model fit, but the resulting effect is small.

Table 3
Goodness of Fit Statistics for the Baseline and Expanded Models

Model	AIC	BIC	logLik
Single fixation, baseline	3140	3918	-1483
Single fixation, expanded	3038	4005	-1411
Gaze fixation, baseline	48270	49081	-24049
Gaze fixation, expanded	47904	48961	-23840

Note. All models contain the same random effects, representing intersubject variability in the effects of relative fixation position, word length, lexical status and predictability. The specific statistics displayed are the Akaike Information Criterion (AIC), the Bayesian Information Criterion (BIC), and the raw log likelihood (logLik), as calculated by the *lmer* program.

Gaze durations

In the baseline model for gaze durations, almost all the lexical features of words n , $n - 1$, and $n + 1$ have a significant effect on the fixation duration on word n . The signs of the main effects of the lexical parameters (linear terms only) are summarised in Table 4.

Table 4
Main Effects of Lexical Parameters on Gaze Durations

Parameter	Word $n-1$	Word n	Word $n+1$
l	↘	↘	↘
g	↘	↘	↘
h	↘	↘	↗
p	↘	↘	↗
x	↘	↗	.
new: q	↘	.	.

Note. Only the effects of the linear terms are shown here. If increasing the parameter increases fixation duration, then a '↗' is displayed, or if the opposite is true, then a '↘' is displayed. If there is no statistically significant effect (at the 99.7% confidence level), then a '.' is shown. For example, when l_n is smaller (meaning that the current word is longer, since l represents reciprocal length), the fixation duration on the current word is longer.

As in the case of single fixation durations, we add the three predictors q_n , q_{n-1} , and q_{n+1} to the baseline model, and then successively add interactions. Again, the addition of random effects to the baseline and expanded models has little influence on the pattern of differences between the models.

The only additional main effect from the new predictors comes from q_{n-1} , which has a negative effect on fixation duration. The inclusion of the new predictors leads to compensatory changes of the frequency predictors g_n and g_{n-1} ; these changes occur because of the high covariances between the frequency and CCP predictors (see Table 5). Apart from these expected changes, there is also a weakening of the main effect of the predictability p_n and an increase in the main effect of the skipping status s_{n-1} .

Table 5
Covariance Matrix for the Single Fixation Durations Model

	g_{n-1}	g_{n+1}	p_n	p_{n-1}	p_{n+1}	x_n	x_{n-1}	x_{n+1}	q_n	q_{n-1}	q_{n+1}
g_n	↘	.	.
g_{n-1}	↘	.
g_{n+1}	↗	.	.	↘
p_n
p_{n-1}
p_{n+1}
x_n
x_{n-1}
x_{n+1}	↘
q_n
q_{n-1}

Note. The predictors shown are content word frequency (g), predictability (p), lexical status (x) and CCP (q), for words n , $n - 1$, and $n + 1$. A positive covariance more positive than 0.25 is represented by a '↗', and a negative covariance more negative than -0.25 is represented by a '↘'. Double arrows signify covariances with magnitudes larger than 0.50. Small covariances (magnitude less than 0.25) are represented by a '.'.

Most of the interactions from the baseline model remain unchanged after the addition of interactions with the new predictors. Notable exceptions are $l_n : g_n$, $l_n : h_n$, $x_{n-1} : p_n$, $g_n : s_{n-1}$, and $g_{n-1} : s_{n-1}$, which change in the expected directions to compensate for the addition of the new interactions. Like in the single fixation case, the inclusion of the new predictors only causes a small improvement in the model fit (see Table 3), which is counterbalanced by the increase in complexity of the model; there is only a small decrease in the Bayesian Information Criterion.

Discussion

Our study has shown that CCP plays a role similar to that of frequency when used as a predictor for fixation durations. In addition, CCP does not seem to affect the role of predictability as a predictor of fixation durations. Addition of CCP to the fixation duration models only marginally improved the fit of the models.

The conclusion that CCP and frequency play similar roles is not trivial, since CCP and frequency are conceptually quite different word properties: CCP is context dependent, while frequency is context independent. Much of the similarity between these two measures can be attributed to the high correlation between them. Why does this high correlation occur?

Remember that CCP for a target word w_j with context C is defined to be $p(w_j|C)$. If word w_j occurs very frequently, then $p(w_j|C) \approx 1$, independent of context. This means that the CCP of highly frequent words must have a small dynamic range. This can be easily seen in Figure 2. Since, by definition, highly frequent words must occur often in a text corpus, there will be many corpus words with high frequency and CCP close to one. The dependence of the dynamic range of CCP on frequency, combined with the large number of words with limited dynamic range, leads to a high correlation between CCP and frequency.

In spite of the high correlation between CCP and frequency, the expanded model fits suggest that CCP does not *only* encode frequency information when used for predicting eye movements. Almost all of the differences in structure between the expanded and baseline models correspond to effects on the current word or the word neighbour to the left. This shows that CCP contributes lag and immediacy effects, but not successor effects, when it is used as a predictor for fixation durations.

Are these results specific to CCP? What about other simple statistical representations of semantic relatedness? One such predictor used by Turney (2001) was an Internet-based *pointwise mutual information* measure (PMI_{ij}), which can be defined in the following way:

$$PMI_{ij} = \log_2 \left(\frac{p(w_i \cap w_j)}{p(w_i)p(w_j)} \right).$$

This can be rewritten as

$$PMI_{ij} = (\log p(w_j|w_i) - \log f(w_i) + \log N) / \log 2,$$

where N is the total number of documents. From this rearrangement, we can see that pointwise mutual infor-

mation is just a linear combination of the measures that we have already included in the regression model, and thus cannot contribute new information. Thus, if another statistical measure of semantic relatedness is to be investigated as a predictor of fixation durations, it needs to have a conceptually different basis, like, for example, Latent Semantic Analysis (Landauer & Dumais, 1997).

It is somewhat disappointing that CCP is so unsuitable to act as a replacement of predictability, especially considering the previous studies suggesting that it might be able to detect simple semantic relationships. However, a task where subjects are required to read for comprehension is far different from one where subjects need to decide on synonymity or judge plausibility. The slight improvement of the gaze duration model after the addition of lag and immediacy CCP effects might hint that subjects only calculate a measure of semantic similarity on multiply fixated words, where there is some need to re-evaluate a word being read, but this would need to be tested more carefully. From the two main results, that (i) the addition of CCP does not significantly improve the fit of the single fixation duration model, and (ii) there are no successor effects containing CCP, it seems that subjects do not have ready access to a simple statistical representation of "semantic similarity" while reading for meaning.

Acknowledgements

This research was performed while James Ong was based at the University of Potsdam in the group of Reinhold Kliegl, and funded by the Deutsche Forschungsgemeinschaft (KL 955/6). The writing of this paper has been supported by a grant (FWF L425-N15) from the Fonds zur Förderung der wissenschaftlichen Forschung (FWF). Thomas Haslwanter provided support and constructive criticism that helped to shape this paper.

References

- Bates, D. (2007). lme4: Linear mixed-effects models using S4 classes [Computer software manual]. (R package version 0.99875-9)
- Boston, M. F., Hale, J., Kliegl, R., Patil, U., & Vasishth, S. (2008). Parsing costs as predictors of reading difficulty: An evaluation using the Potsdam Sentence Corpus. *Journal of Eye Movement Research*.

- Frisson, S., Rayner, K., & Pickering, M. J. (2005). Effects of contextual predictability and transitional probability on eye movements during reading. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 31(5), 862–877.
- Keller, F., & Lapata, M. (2003). Using the web to obtain frequencies for unseen bigrams. *Computational Linguistics*, 29(3), 459–484.
- Kennedy, A., & Murray, W. S. (1987). Spatial coordinates and reading: Comments on Monk (1985). *The Quarterly Journal of Experimental Psychology Section A*, 39(4), 649–656.
- Kliegl, R. (2007). Toward a perceptual-span theory of distributed processing in reading: A reply to Rayner, Pollatsek, Drieghe, Slattery, and Reichle (2007). *Journal of Experimental Psychology: General*, 136(3), 530–537.
- Kliegl, R., Grabner, E., Rolfs, M., & Engbert, R. (2004). Length, frequency, and predictability effects of words on eye movements in reading. *European Journal of Cognitive Psychology*, 16(1/2), 262–284.
- Kliegl, R., Nuthmann, A., & Engbert, R. (2006). Tracking the mind during reading: The influence of past, present, and future words on fixation durations. *Journal of Experimental Psychology: General*, 135(1), 12–35.
- Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato's problem: The Latent Semantic Analysis theory of the acquisition, induction, and representation of knowledge. *Psychological Review*, 104, 211–240.
- McDonald, S. A., & Shillcock, R. C. (2003). Eye movements reveal the on-line computation of lexical probabilities during reading. *Psychological Science*, 14, 648–652.
- Pynte, J., & Kennedy, A. (2006). An influence over eye movements in reading exerted from beyond the level of the word: Evidence from reading English and French. *Vision Research*, 46, 3786–3801.
- Pynte, J., New, B., & Kennedy, A. (in press). A multiple regression analysis of syntactic and semantic influences in reading normal text. *Journal of Eye Movement Research*.
- R Development Core Team. (2008). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. (ISBN 3-900051-07-0, <http://www.R-project.org/>)
- Rayner, K., Ashby, J., Pollatsek, A., & Reichle, E. D. (2004). The effects of frequency and predictability on eye fixations in reading: Implications for the E-Z reader model. *Journal of Experimental Psychology: Human Perception and Performance*, 30, 720–732.
- Taylor, W. L. (1953). Cloze procedure: A new tool for measuring readability. *Journalism Quarterly*, 30, 415–433.
- Turney, P. D. (2001). Mining the web for synonyms: PMI-IR versus LSA on TOEFL. In L. D. Raedt & P. A. Flach (Eds.), *Machine learning: ECML 2001, 12th European Conference on Machine Learning, Freiburg, Germany, September 5-7, 2001, Proceedings* (Vol. 2167, pp. 491–502). Springer.
- White, S. J. (2008). Eye movement control during reading: Effects of word frequency and orthographic familiarity. *Journal of Experimental Psychology: Human Perception and Performance*, 34(1), 205–223.
- White, S. J., & Liversedge, S. P. (2006). Linguistic and nonlinguistic influences on the eyes' landing positions during reading. *The Quarterly Journal of Experimental Psychology*, 59(4), 760–782.